

ARE PRIVATE SCHOOLS BETTER THAN PUBLIC SCHOOLS? APPRAISAL FOR IRELAND BY METHODS FOR OBSERVATIONAL STUDIES¹

BY DANNY PFEFFERMANN AND VICTORIA LANDSMAN

*Hebrew University of Jerusalem and University of Southampton,
 and National Cancer Institute*

In observational studies the assignment of units to treatments is not under control. Consequently, the estimation and comparison of treatment effects based on the empirical distribution of the responses can be biased since the units exposed to the various treatments could differ in important unknown pretreatment characteristics, which are related to the response. An important example studied in this article is the question of whether private schools offer better quality of education than public schools. In order to address this question, we use data collected in the year 2000 by OECD for the *Programme for International Student Assessment* (PISA). Focusing for illustration on scores in mathematics of 15-year-old pupils in Ireland, we find that the raw average score of pupils in private schools is higher than of pupils in public schools. However, application of a newly proposed method for observational studies suggests that the less able pupils tend to enroll in public schools, such that their lower scores are not necessarily an indication of bad quality of the public schools. Indeed, when comparing the average score in the two types of schools after adjusting for the enrollment effects, we find quite surprisingly that public schools perform better on average. This outcome is supported by the methods of instrumental variables and latent variables, commonly used by econometricians for analyzing and evaluating social programs.

1. Introduction. In observational studies the assignment of units to treatments often depends on latent variables that are related to the response variable even when conditioning on known covariates. Consequently, a direct comparison of the response distributions (given the model covariates) or moments of these distributions between treatment groups could be bi-

Received August 2010; revised December 2010.

¹Supported by the Israel Science Foundation Grant 1277/05.

Key words and phrases. Average treatment effect, goodness of fit, identifiability, instrumental variables, private-dependent schools, propensity scores, sample distribution.

This is an electronic reprint of the original article published by the
 Institute of Mathematical Statistics in *The Annals of Applied Statistics*,
 2011, Vol. 5, No. 3, 1726–1751. This reprint differs from the original in pagination
 and typographic detail.

ased and result in wrong conclusions. An important example studied in this article is the question of whether private schools offer better quality of education than public schools. This question has important impact on educational policy and public finance [Hanushek (2002)]. It is known that pupils enrolling to the two types of schools differ in their family background and other characteristics related to their scholastic achievements, such that a raw comparison of the scores of pupils attending the two types of schools can be misleading. In an attempt to deal with this question, we use data collected in the year 2000 by OECD for the *Programme for International Student Assessment* (PISA). The purpose of this program is to study and compare the proficiency of pupils aged 15 from over more than 30 countries in mathematics, science and reading. In this article we focus on scores in mathematics in Ireland and estimate the difference in the average score between the two types of schools by adjusting for the quality of pupils enrolling to them and for the effects of known covariates.

We start by applying several existing methods for observational studies to the data, which are described in Section 3, and find, similarly to Vanderberghe and Robin (2004), that some of these methods produce estimates of different magnitude and sign. We attempt to resolve this conflict by developing and applying a new method for inference from observational data, which extends recent methodology for analyzing sample survey data. The method derives the *sample distribution* of the observed response under a given treatment (score in mathematics in a given type of school in our application) as a function of the distribution that would be obtained under a strongly ignorable assignment of subjects to treatments (assumptions $SI(a)$, $SI(b)$ in Section 3), and the assignment probability, which is allowed to depend on the response value. The use of this approach is established by showing that the sample distribution is identifiable under some general conditions. The goodness of fit of the sample distribution can be tested by standard test statistics since it refers to the observed data.

By fitting the sample distribution to the observed data, we can estimate the distribution under strongly ignorable assignment to treatments, and the assignment probabilities, which are then used for estimating population means or contrasts between them. Our approach permits also testing some of the assumptions underlying other methods for analyzing observational data, thus enabling us to understand better why different methods yield different answers in our application.

Section 2 describes the PISA data and defines the problem underlying this study more formally. Section 3 overviews some of the existing methods for observational studies and shows the results obtained when applying them to the PISA data for Ireland. We also consider “probability weighted” versions of the estimators, which account for unequal sample selection probabilities that are possibly related to the response and may thus bias the inference. Computation of these estimators yields very similar estimates to

the estimates obtained under the standard methods. Section 4 presents the proposed approach and shows the results obtained when applying it to the PISA data. The main outcome of this analysis is that after controlling for the effect of the enrollment process, the public schools actually outperform the private schools in the average math score, suggesting better quality of education. Here also we extend the method to account for unequal selection probabilities and obtain similar estimates. Section 5 overviews the main theoretical properties of the new approach. The technical derivations are presented in Supplements C and D of the supplementary material [Pfeffermann and Landsman (2011)]. We conclude with a brief summary in Section 6.

2. Data used for application and formulation of the problem.

2.1. Sampling design and response values. In order to compare the private and public schools, we use data collected in *Ireland* in the year 2000 by OECD for PISA.

Sampling design. PISA uses in most countries a stratified two-stage sampling design. The strata are defined by the size of the school, type of school and gender composition. In each stratum, a probability proportional to size (PPS) sample of schools was selected with the size defined by the number of 15-year-old pupils enrolled in the school. A minimum of 150 schools has been selected in each country, or all the schools if there are less. In the second stage an equal probability sample of 35 pupils from the corresponding age group was drawn from each of the sampled schools (or all the pupils in schools with less than 35 pupils aged 15). By this sampling design, pupils included in the sample do not have equal selection probabilities and each pupil is assigned therefore a sampling weight. The weight is the reciprocal of the pupil's sample inclusion probability, adjusted for nonparticipation of schools and nonresponse of pupils.

PISA distinguishes between two types of private schools: private-dependent schools where the government contributes 50% or more to the school core funding and private-independent schools with less than 50% government funding. The sample from Ireland consists of 54 public schools, 79 private-dependent schools and only 4 private-independent schools and, hence, in this paper we do not distinguish between the two types of schools and refer to them simply as private schools. For more information on the PISA sampling design and weighting, see Adams and Wu (2002).

Computation of response values. The pupils' proficiencies (scores in mathematics in our case) are not observed directly in the PISA study and are viewed as missing data, which are imputed from the item responses $d_j = (d_{1j}, d_{2j}, \dots, d_{mj})$, where $d_{ij} = 1$ if pupil j answers correctly question i of the examination and $d_{ij} = 0$ otherwise, $i = 1, \dots, m$. PISA uses two approaches for imputing the scores: a maximum likelihood approach and a multiple

imputation approach. In this paper we used the imputed values obtained under the second approach. See Appendix A for the imputation model. The PISA database contains five sets of imputed values. We standardized the imputed scores in each set by dividing them by their empirical standard deviation and then defined the response value to be the average of the five standardized values. After standardization and averaging, the range of the response values is approximately from 1 to 10. We compared the use of the average values to the results obtained when analyzing each of the five sets of standardized values separately and then combining the results using multiple imputation theory and obtained very similar results in all the analyses performed. Consequently, in this paper we restrict to the average response since it is convenient to have a single working model when simulating new observations, which is needed for the goodness-of-fit tests discussed later.

2.2. Formulation of the problem. The formulation of the problem for the PISA data follows what is known in the literature as the *counterfactual approach*. By this approach, every unit in the population is potentially exposed to every treatment. See, for example, Rubin (1974), Rosenbaum and Rubin (1983), Smith and Sugden (1988) and Rosenbaum (2002).

Let U define the population of 15-year-old pupils in Ireland. Every pupil $i \in U$ has two *potential* responses: Y_{1i} —the proficiency score if the student attends a private school, and Y_{0i} —the proficiency score if the student attends a public school. Let \mathbf{x} denote a set of k known covariates (background characteristics) that affect the responses, with values \mathbf{x}_i for pupil i . The (potential) population mean score in private schools (hereafter the treatment group) is defined as $\mu_1 = \frac{1}{N} \sum_{i=1}^N E_p[Y_{1i}|\mathbf{x}_i]$, where N is the population size and the expectation $E_p(\cdot)$ is with respect to the population model holding for the responses. The population mean score in public schools (hereafter the control group) is defined accordingly as $\mu_0 = \frac{1}{N} \sum_{i=1}^N E_p[Y_{0i}|\mathbf{x}_i]$. In many observational studies, contrasts between the parameters μ_1 and μ_0 are of primary interest. In this paper we focus on estimating the difference between the mean score in private and public schools, defined as

$$(2.1) \quad \tau = \mu_1 - \mu_0 = \frac{1}{N} \sum_{i=1}^N E_p[Y_{1i}|\mathbf{x}_i] - \frac{1}{N} \sum_{i=1}^N E_p[Y_{0i}|\mathbf{x}_i].$$

The contrast τ is known in the literature as the *average treatment effect* (ATE).

In practice, every unit in the population is only exposed to one treatment. Also, it is rarely the case that all the population units participate in the study. The observed data refer therefore to a sample S of size n , which in our application is divided into the two subsamples S_1 and S_0 , where S_1 (S_0) is the subsample of pupils attending private (public) schools. For every pupil $i \in S$ we observe therefore y_{1i} if $i \in S_1$ or y_{0i} if $i \in S_0$.

Denote by π_i the probability that pupil i is included in the sample S and by \tilde{p}_{ti} the probability that *sampled* pupil i is enrolled in school of type t . The sample inclusion probabilities π_i (or the sampling weights $w_i = 1/\pi_i$, with possible adjustments for nonresponse or calibration) are typically known for the sampled units, as is the case in the PISA survey, but the treatment assignment probabilities, \tilde{p}_{ti} , are usually unknown and may depend on latent variables that are related to the response, Y_{ti} . As is well known and illustrated later, if the effect of these latent variables on the response is not accounted for by the observed covariates, the resulting estimators of the population parameters can be highly biased.

REMARK 1. The sample inclusion probabilities may likewise be related to the response values and thus bias the inference if not accounted for adequately. This is known in the survey sampling literature as *informative sampling*. Smith and Sugden (1988) define conditions on the sampling design and the treatment assignment process that warrant ignoring them in the inference process. As shown in subsequent sections, there is no evidence for informative sampling with the kind of models and inference methods applied to the PISA data from Ireland.

3. Existing methods, application to PISA data. In what follows we focus on the estimation of the ATE defined by (2.1), assuming that the sample selection probabilities are not related to the response variable and the covariates, and hence that there are no sampling effects. This is the common assumption in the literature even though seldom stated explicitly. After describing several methods in common use and applying them to the PISA data, we show the results obtained when extending the methods to the case where the sample is selected with known unequal probabilities that might be related to the response and/or the covariates and compare the results with the results obtained when ignoring the sample selection. Let T define the indicator of the treatment group ($T = 1$ for private schools, $T = 0$ for public schools). Rosenbaum and Rubin (1983) establish two conditions that warrant ignoring the treatment assignment in the inference process when conditioning on \mathbf{x} :

- SI(a): the assignment T and the response values (Y_1, Y_0) are independent given the covariates, \mathbf{x} , for every unit (pupil),
- SI(b): $0 < \Pr(T = 1|\mathbf{x}) < 1$ for every possible \mathbf{x} .

Conditions SI(a) and SI(b) define a *strongly ignorable assignment process* given the covariates. When the assignment is strongly ignorable, it permits the application of a number of simple estimation techniques, which we review in Section 3.1. In Sections 3.2 and 3.3 we consider the latent variables method (LV) and the use of instrumental variables (IV), which do not assume strong ignorability assumptions.

3.1. Methods for strongly ignorable treatment assignments.

Regression estimator. Suppose that the true relationship between Y and \mathbf{x} in the population has the general form $Y_t = r_t(\mathbf{x}) + u_t$, $E_p(u_t|\mathbf{x}) = 0$ for some functions $r_t(\mathbf{x})$, $t = 0, 1$, where $E_p(\cdot)$ is the expectation under a strongly ignorable assignment. Then, the ATE is $(\bar{r}_1 - \bar{r}_0)$, where $\bar{r}_t = \frac{1}{N} \sum_{i=1}^N r_t(\mathbf{x}_i)$. When the expectations $r_t(\mathbf{x})$ are linear, $r_t(\mathbf{x}) = \mathbf{x}'\beta_t$, then under the assumptions SI(a), SI(b) the regression coefficients can be estimated by ordinary least squares (OLS) and the ATE estimator takes the form

$$(3.1) \quad \hat{\tau}_{\text{OLS}} = \bar{\mathbf{x}}'(\hat{\beta}_{1,\text{OLS}} - \hat{\beta}_{0,\text{OLS}}),$$

where $\bar{\mathbf{x}}' = (\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_K) = \sum_{i=1}^n \mathbf{x}_i/n$ and $\hat{\beta}_{t,\text{OLS}}$ is the OLS estimator in subsample S_t .

Matching estimator. Another procedure in common use is to match the units from the treatment and control groups based on the covariates and then compare the responses. Matching procedures are widely discussed in the literature; see, for example, Rosenbaum (2002). They do not require specifying the form of the functions $r_t(\mathbf{x})$. Abadie and Imbens (2006) consider the following matching estimate with replacement. Denote by $J_{iM}(t)$ the indices of the M closest matches in S_t for unit $i \in S_{1-t}$, $t = 0, 1$. Define for unit $i \in S_{1-t}$, $\hat{y}_{(1-t),i} = y_{(1-t),i}$ and $\hat{y}_{ti} = \frac{1}{M} \sum_{j \in J_{iM}(t)} y_{tj}$. Estimate

$$(3.2) \quad \hat{\tau}_M = \frac{1}{n} \sum_{i=1}^n (\hat{y}_{1i} - \hat{y}_{0i}).$$

Other methods use probability weighting with the weights defined by the inverse of the “propensity score,” $e(\mathbf{x}) = \Pr(T = 1|\mathbf{x})$. Rosenbaum and Rubin (1983) show that the conditions SI(a), SI(b) for strong ignorability imply the same conditions when \mathbf{x} is replaced by $e(\mathbf{x})$, thus validating the use of propensity scores for ATE estimation. In practice, the propensity scores are unknown and are estimated by fitting logistic or probit models, or by use of nonparametric techniques [McCaffrey, Ridgeway and Morral (2004)]. Below we describe two ATE estimators that use the estimated propensity scores, $\hat{e}(\mathbf{x}_i)$, for weighting.

Brewer–Hajek (B–H) estimator. This estimator resembles the familiar Brewer–Hajek [Brewer (1963); Hajek (1971)] estimator in survey sampling. Let $T_i = 1(0)$ if unit $i \in S_1$ ($i \in S_0$) and define $Y_i = T_i Y_{1i} + (1 - T_i) Y_{0i}$. The B–H estimator is

$$(3.3) \quad \hat{\tau}_{\text{B-H}} = \left(\sum_{i=1}^n \frac{T_i}{\hat{e}(\mathbf{x}_i)} \right)^{-1} \sum_{i=1}^n \frac{T_i Y_i}{\hat{e}(\mathbf{x}_i)} - \left(\sum_{i=1}^n \frac{1 - T_i}{1 - \hat{e}(\mathbf{x}_i)} \right)^{-1} \sum_{i=1}^n \frac{(1 - T_i) Y_i}{1 - \hat{e}(\mathbf{x}_i)}.$$

Doubly-robust (DR) estimator. If the population expectation $E_p(Y_t|\mathbf{x})$ can be modeled by some function $r_t(\mathbf{x})$, then by SI(a), $r_t(\mathbf{x})$ is also the sample

expectation and the ATE can be estimated as

$$(3.4) \quad \hat{\tau}_{\text{DR}} = \frac{1}{n} \sum_{i=1}^n \frac{T_i Y_i - [T_i - \hat{e}(\mathbf{x}_i)] \hat{r}_1(\mathbf{x}_i)}{\hat{e}(\mathbf{x}_i)} - \frac{1}{n} \sum_{i=1}^n \frac{(1 - T_i) Y_i + [T_i - \hat{e}(\mathbf{x}_i)] \hat{r}_0(\mathbf{x}_i)}{1 - \hat{e}(\mathbf{x}_i)}.$$

The estimator (3.4) has the “double-robustness” property of being consistent even if only the model assumed for the propensity scores or the population expectations are correctly specified [Lunceford and Davidian (2004)]. Qin and Zhang (2007) consider another estimator that has a somewhat stronger robustness property.

3.2. Latent variable models. This method specifies the joint distribution of the outcome and the treatment selection by use of latent variable (LV) models. The model assumes the following:

- LV(a)—a structural equation for the population outcomes of the form, $Y_t = r_t(\mathbf{x}) + u_t$; $E_p(u_t|\mathbf{x}) = 0$, $t = 0, 1$, and
- LV(b)—a latent variable u_T and an assignment rule T satisfying $T = 1[l(\mathbf{v}; \alpha) + u_T > 0]$, where $1[\cdot]$ is the indicator function and $l(\mathbf{v}; \alpha)$ is a given function of known covariates \mathbf{v} , governed by a vector parameter α .

The covariates in \mathbf{v} may include some of the covariates in \mathbf{x} , but it is generally recommended that $\mathbf{v} \neq \mathbf{x}$ to avoid colinearity problems in the estimation process; see below. The random variables (u_1, u_0, u_T) are dependent. Under these assumptions, $E_{S_t}(Y_t|\mathbf{x}) = E(Y_t|\mathbf{x}, T = t) = r_t(\mathbf{x}) + E(u_t|\mathbf{x}, T = t) \neq r_t(\mathbf{x}) = E_p(Y_t|\mathbf{x})$, since $E(u_t|\mathbf{x}, T = t) \neq 0$. However, assuming that $r_t(\mathbf{x}) = \mathbf{x}'\beta_t$, (u_1, u_0, u_T) are jointly normal and $l(\mathbf{v}; \alpha) = \mathbf{v}'\alpha$, β_t can be estimated by the two-stage Heckman’s method [Maddala (1983)], yielding

$$(3.5) \quad \hat{\tau}_{\text{LV}} = \bar{\mathbf{x}}'(\hat{\beta}_{1,\text{LV}} - \hat{\beta}_{0,\text{LV}}),$$

where $\hat{\beta}_{t,\text{LV}}$ is the LV estimator of β_t , $t = 0, 1$. Heckman and Vytlacil [(2006), Ch. 70] refer to the latter model as the Generalized Roy Model and discuss semi-parametric econometric models, which relax some of the assumptions of this model.

3.3. Instrumental variables models. Let $Y_t = \mathbf{x}'\beta_t + u_t$; $E_p(u_t|\mathbf{x}) = 0$, $t = 0, 1$. Then for unit $i \in S$,

$$(3.6) \quad Y_i = T_i \mathbf{x}'_i \beta_1 + (1 - T_i) \mathbf{x}'_i \beta_0 + u_i = \tilde{\mathbf{x}}'_i \theta + u_i,$$

where $Y_i = T_i Y_{1i} + (1 - T_i) Y_{0i}$, $\tilde{\mathbf{x}}'_i = (T_i \mathbf{x}'_i, (1 - T_i) \mathbf{x}'_i)$, $\theta = (\beta'_1, \beta'_0)'$ and $u_i = T_i u_{1i} + (1 - T_i) u_{0i} = u_{0i} + T_i(u_{1i} - u_{0i})$. In observational studies u_1 and u_0 are

correlated with T and, hence, θ cannot be estimated consistently from (3.6) without additional assumptions. Below we define a set of plausible assumptions warranting that the ATE is estimated consistently. See Wooldridge (2002) for discussion of this and alternative sets of assumptions. Assume the availability of an instrument h satisfying the following:

- IV(a)— $E_p(Y_t|x, h) = E_p(Y_t|x)$ (the population expectation under strongly ignorable assignment does not depend on the instrument, given the covariates);
- IV(b)— $E_p[T(u_1 - u_0)|x, h] = 0$ (the assignment and the counterfactual gain in the error terms are uncorrelated given the covariates and the instrument);
- IV(c)— $\Pr(T = 1|x, h) = g(x, h) \neq \Pr(T = 1|x)$ (the assignment probabilities depend on the instrument and possibly on x).

Multiplying both sides of (3.6) by the column vector $z_i = (g_i x'_i, (1 - g_i)x'_i)'$, where $g_i = g(x_i, h_i)$ and taking expectations yields $E_p(z_i Y_i | x_i, h_i) = E_p(z_i \tilde{x}'_i | x_i, h_i) \theta$, since under the model and IV(b), $E_p(z_i u_i | x_i, h_i) = 0$. The IV estimator of θ computed from all the observations is $\hat{\theta}_{IV} = (\sum_{i=1}^n \hat{z}_i \tilde{x}'_i)^{-1} \sum_{i=1}^n \hat{z}_i Y_i$, where $\hat{z}_i = (\hat{g}_i x'_i, (1 - \hat{g}_i)x'_i)'$. The estimator $\hat{g}_i = \hat{g}(x_i, h_i)$ is commonly obtained by fitting probit or logit models. The ATE estimator is

$$(3.7) \quad \hat{\tau}_{IV} = \bar{x}'(\hat{\beta}_{1,IV} - \hat{\beta}_{0,IV}),$$

with $\hat{\beta}_{t,IV}$ defined by $\hat{\theta}_{IV}, t = 0, 1$.

REMARK 2. Condition IV(c) is testable from the data, but conditions IV(a) and IV(b) relate to unobservable quantities and cannot in general be tested directly. Imbens and Angrist (1994) show that for a binary instrument, if condition IV(b) is not satisfied, then under a weaker *monotonicity condition*, $\hat{\tau}_{IV}$ estimates the treatment effect for a subpopulation consisting of units for which the treatment status would be altered by the instrument. This treatment effect is called *local average treatment effect* (LATE).

3.4. *Application of the methods to PISA data for Ireland.* We applied the methods reviewed so far to the PISA data for Ireland described in Section 2. The sample consists of 1,244 pupils from private schools and 694 pupils from public schools. Six covariates were found to be significant in at least one of the models described in Section 4: gender (GEN; 1 for girls, 0 for boys), mother's education (ME; 1 for high education, 0 otherwise), family socio-economic index (SEI), index of home educational resources (HER), average socio-economic index of the pupil's schoolmates [SES; proposed by Vandenberghe and Robin (2004) to account for potential peer effects], and school location (S.loc; 1 if the school is located in an urban area, 0 otherwise). The continuous covariates have been standardized. To warrant fair

TABLE 1
Estimates of ATE and standard errors under existing methods

Method	$\hat{\tau}_D$	$\hat{\tau}_{OLS}$	$\hat{\tau}_M$	$\hat{\tau}_{B-H}$	$\hat{\tau}_{DR}$	$\hat{\tau}_{LV}$	$\hat{\tau}_{IV}$
Estimate	0.36	0.12	0.21	0.16	0.17	-0.49	-0.61
Std. error	0.05	0.05	0.05	0.05	0.05	0.19	0.24

comparability between the various methods, we included for the first four methods [equations (3.1)–(3.4)] all the six covariates in both the regressions and the models used for computing the propensity scores. For the LV and IV methods we included all the covariates except for S.loc in the regressions and all the covariates including S.loc in the school selection models (see Remark 3). Vandenberghe and Robin (2004) considered additional covariates, but these were not found to be significant in our analysis.

REMARK 3. The variable *school location* was used by Vandenberghe and Robin (2004) as an instrumental variable. The authors show that it has a significant effect on the probability of attending private schools, thus satisfying the condition IV(c) in Section 3.3. However, the approaches considered in the literature for observational studies do not permit testing that the school location is exogenous to the pupil’s proficiency given the other covariates, as required by condition IV(a), because this condition refers to the population models of the unobservable potential responses. The authors argue that this requirement is plausible, using similar arguments to Hoxby (2000). See Section 4.6 for how we can test this condition under the approach proposed in Section 4.

Table 1 presents the ATE estimates and their standard errors. The first estimate, $\hat{\tau}_D = \bar{y}_1 - \bar{y}_0$, is the crude difference between the simple sample means in the two types of schools. The matching estimator is computed based on $M = 4$ matches. We considered several matching estimates as obtained under different metrics for finding the matches, with and without adjustments for imperfect matching, and obtained very close results in all the cases. For the instrumental variables method we used the school location as the instrument.

The estimates $\hat{\tau}_M$, $\hat{\tau}_{LV}$ and $\hat{\tau}_{IV}$ were computed by using the functions *nnmatch*, *treatreg* and *ivreg* of the Stata software [StataCorp (2004)]. The remaining estimates were programmed using the R software [R Development Core Team (2004)]. Estimation of the standard errors of the matching estimators and the LV and IV estimators is incorporated in the Stata functions. See Abadie and Imbens (2006) and Wooldridge (2002) for details. Estimation of the standard errors of the Brewer–Hajek estimator and the doubly robust estimator is developed by Lunceford and Davidian (2004). The esti-

mated standard errors account for the error distributions of the responses under the respective models.

The first notable outcome in Table 1 is that the difference $\hat{\tau}_D$ between the simple sample means in the two types of schools is positive, which we anticipated because the more able pupils tend to enroll in private schools. The next four methods from left, which assume strongly ignorable assignment given the covariates, likewise produce small positive ATE estimates. By contrast, the IV and LV methods, which account for treatment assignment effects not explained by the covariates, produce negative estimates, with much larger absolute values, suggesting that the public schools actually perform better after accounting for the school selection effects. A similar outcome is obtained under the approach proposed in Section 4. The use of this approach explains also why the LV and IV methods are more appropriate for this data.

REMARK 4. Vandenberghe and Robin (2004) computed what is known in the econometric literature as “the average treatment effect for the treated (ATT),” using the same data and some of the methods reviewed before, and obtained similar results to the results in Table 1. Dronkers and Avram (2010) computed the ATT for reading scores using PISA data for all the countries by applying several variants of propensity scores matching. The ATT estimates for Ireland in this study are positive, same as the ATE estimates for the scores in Mathematics based on propensity scores presented in Table 1 ($\hat{\tau}_{B-H}$ and $\hat{\tau}_{DR}$).

3.5. *Probability weighted estimators for PISA data.* So far we ignored the sample selection process when computing the estimates in Section 3.4. The question arising is whether this is justified in the present study. We emphasize again that if the distribution of the response in the sample is affected by the sample selection scheme, the sampling is informative and failing to account for the sampling effects may bias the inference. In fact, even if only the distribution of the covariates in the sample is different from their population distribution, some of the ATE estimators may already be biased. Pfeiffermann and Sverchkov (2009) review several existing approaches to account for possible sampling effects in the inference process. In this study we applied what is known as probability weighting, which basically consists of inflating each sample observation proportionally to its sampling weight. The idea of probability weighting is to obtain estimators that are consistent under the randomization (repeated sampling) distribution for the corresponding “census estimates” that would be computed if all the population values had been observed. The census estimates are free of sampling effects.

We computed the probability weighted estimators (PWE) for all the methods considered so far. See Supplement A in the supplementary material [Pfeiffermann and Landsman (2011)] for the derivation of these estimators. As

TABLE 2
Unweighted and weighted estimators of schools mean score and ATE

Method	Private schools		Public schools		ATE	
	UNWEI	WEI	UNWEI	WEI	UNWEI	WEI
Simple difference	6.28	6.29	5.92	5.92	0.36	0.37
Regression	6.21	6.26	6.09	6.12	0.12	0.14
Matching	6.25	6.26	6.04	6.03	0.21	0.23
Brewer–Hajek (B–H)	6.24	6.26	6.08	6.06	0.16	0.20
Doubly robust (DR)	6.23	6.25	6.06	6.07	0.17	0.18
Instrumental variable	6.00	6.02	6.61	6.52	−0.61	−0.50
Latent variable	6.00	6.02	6.49	6.41	−0.49	−0.39

a first step we computed the unweighted and probability weighted estimators (in parenthesis) of the population means of the covariates and obtained the following results (the covariates are defined in Section 3.4): GEN: 0.53 (0.52), ME: 0.61 (0.61), SEI: 0.00 (−0.016), SES: 0.00 (−0.016), HER: 0.00 (0.002), S.loc: 0.40 (0.39). As can be seen, the two sets of estimators are very close.

Table 2 shows for each of the methods the unweighted (UNWEI) and probability weighted (WEI) estimators of the mean score in the private and public schools, and the corresponding ATE estimator. The results in Table 2 indicate that the PWE of the mean score as obtained under the various methods are very similar to the corresponding unweighted estimators. This is definitely true for the private schools, but even for the public schools the largest difference between the weighted and unweighted estimate is less than 2%. The very small differences between the weighted and unweighted estimates in each type of school translate into somewhat larger differences in the estimates of the ATE, but not to an extent that affects the inference. Notice in this regard that when computing the conventional 95% confidence intervals for the true ATE based on the unweighted ATE estimates, all the intervals contain the corresponding weighted estimates. In fact, this would be the case even for confidence intervals with confidence level as low as 68%. Our general conclusion from Table 2 is therefore that the sampling process can be ignored when analyzing the PISA data from Ireland by use of the methods considered so far.

4. An alternative approach for observational studies. In this section we propose an alternative approach for ATE estimation, which, as illustrated in Section 4.6, allows also testing the appropriateness of candidate instrumental variables or the use of propensity scores under the assumed model. The approach resembles the LV approach in the sense that it assumes a population model and a model for the treatment selection and applies a combined likelihood resulting from the two models, but all the subsequent develop-

ments are very different. As with the IV and LV methods, the use of this approach does not require strong ignorability assumptions. In what follows we describe the method and apply it to the PISA data assuming noninformative sampling, but later we also consider probability weighted estimation. As before, we consider the case of two groups, $t = 0, 1$.

4.1. The sample distribution. Denote by $f_p(y_{ti}|x_i)$ the *population pdf* for units in treatment group t under a strongly ignorable assignment process. We allow the assignment process to depend on known covariates v , some or all of which may be included in x . Denoting $z = x \cup v$, we assume $f_p(y_{ti}|z_i) = f_p(y_{ti}|x_i)$ and $\Pr(T_i = t|y_{ti}, z_i, i \in S) = \Pr(T_i = t|y_{ti}, v_i, i \in S)$. The *sample pdf* for unit i exposed to treatment t , given the covariates z_i , is obtained by Bayes theorem as

$$(4.1) \quad f_{S_t}(y_{ti}|z_i) = f(y_{ti}|z_i; T_i = t) = \frac{\Pr(T_i = t|y_{ti}, v_i, i \in S) f_p(y_{ti}|x_i)}{\Pr(T_i = t|z_i, i \in S)},$$

where $\Pr(T_i = t|z_i, i \in S) = \int \Pr(T_i = t|y_{ti}, v_i, i \in S) f_p(y_{ti}|x_i) dy_{ti}$.

REMARK 5. It follows from (4.1) that the sample *pdf* is generally different from the population *pdf*, unless $\Pr(T_i = t|y_{ti}, v_i, i \in S) = \Pr(T_i = t|z_i, i \in S)$, in which case the assignment to treatments can be ignored for inference when conditioning on z .

REMARK 6. The probabilities $\Pr(T_i = t|z_i, i \in S)$ are *propensity scores*.

The sample *pdf* defined by (4.1) was shown in recent years to provide a valuable modeling approach for inference from complex sample surveys; see Pfeiffermann and Sverchkov (2009) for review of studies that utilize the sample *pdf* for inference on generalized linear models, testing of distribution functions and prediction of finite population and small area means. The obvious distinction between survey sampling and observational studies is that in survey sampling the sample inclusion probabilities $\pi_i = \Pr(i \in S)$ are usually known, which enables estimating the probabilities $\Pr(i \in S|y_i, v_i)$ and testing the informativeness of the sampling process [Pfeiffermann and Sverchkov (2003, 2009)]. This is generally not the case in observational studies, requiring therefore modeling the parametric form of the probabilities $p_{ti} = \Pr(T_i = t|y_{ti}, v_i, i \in S)$ in (4.1). As discussed below, modeling the sample *pdf* (4.1) allows estimating the unknown parameters governing the *pdf* $f_p(y_{ti}|x_i)$ and the probabilities p_{ti} , and using them for estimating the ATE.

4.2. Estimating the parameters of the sample distribution. So far we suppressed for convenience in the notation the parameters governing the sample *pdf*. Adding the unknown parameters to the notation and assuming that the inclusion in the sample and the assignment to treatments are independent

between units, and that the responses are likewise independent, the sample likelihood for treatment t , based on the sample S_t of size n_t , takes the form

$$(4.2) \quad L_{S_t}[\alpha_t, \theta_t; \{y_{ti}, z_i\}] = \prod_{i=1}^{n_t} \frac{\Pr(T_i = t | y_{ti}, v_i, i \in S; \alpha_t) f_p(y_{ti} | x_i; \theta_t)}{\Pr(T_i = t | z_i, i \in S; \alpha_t, \theta_t)}.$$

Alternatively, the likelihood (4.2) can be replaced by the joint (“full”) likelihood of the treatment selection and the response measurements defined as

$$(4.3) \quad \begin{aligned} L_S[\alpha_t, \theta_t; \{y_{ti}, i \in S_t; z_j, j \in S\}] \\ = \prod_{i=1}^{n_t} \Pr(T_i = t | y_{ti}, v_i, i \in S; \alpha_t) f_p(y_{ti} | x_i; \theta_t) \\ \times \prod_{\substack{j \in S \\ j \notin S_t}} [1 - \Pr(T_j = t | z_j, j \in S; \alpha_t, \theta_t)]. \end{aligned}$$

The likelihood (4.3) has the advantage of comprising the unconditional treatment assignment probabilities for units outside the sample S_t , thus using more information for estimating the model parameters. The full likelihood is often applied for handling informative nonresponse; see, for example, Greenlees, Reece and Zieschang (1982), Rotnitzky and Robins (1997), Gelman et al. [(2004), Chapter 7] and Little (2004).

Replacing the unknown model parameters by their maximum likelihood estimates (*mle*) yields the estimates

$$(4.4) \quad \hat{f}_p(y_{ti} | x_i) = f_p(y_{ti} | x_i; \hat{\theta}_t); \quad \hat{p}_{ti} = \Pr(T_i = t | y_{ti}, v_i, i \in S; \hat{\alpha}_t).$$

4.3. Estimation of population means. When the covariates x_i are known for every unit $i \in U$, then by (4.4), the population means, $\mu_t = \frac{1}{N} \sum_{i=1}^N E_p(Y_{ti} | x_i)$, $t = 0, 1$ (and, consequently, the ATE, $\tau = \mu_1 - \mu_0$) can be estimated by

$$(4.5) \quad \hat{\mu}_{t,p} = \frac{1}{N} \sum_{i=1}^N \hat{E}_p(Y_{ti} | x_i; \theta_t) = \frac{1}{N} \sum_{i=1}^N E_p(Y_{ti} | x_i; \hat{\theta}_t).$$

Note that if $E_p(Y_{ti} | x_i; \theta_t)$ is linear, the computation of (4.5) only requires knowledge of the population means $\bar{X} = (\bar{X}_1, \dots, \bar{X}_k)'$.

When the covariates are unknown for units outside the sample, or the expectation is not linear, then as long as the sampling design is noninformative with respect to the distribution of the covariates, one can use the estimator,

$$(4.6) \quad \hat{\mu}_{t,S} = \frac{1}{n} \sum_{i \in S} \hat{E}_p(Y_{ti} | x_i; \theta_t).$$

Alternatively, one can use in this case the “combined” estimator,

$$(4.7) \quad \hat{\mu}_{t,C} = \frac{\sum_{i=1}^n T_i [Y_{ti} - \hat{E}_p(Y_{ti} | \mathbf{x}_i; \theta_t)] / \hat{p}_{ti}}{\sum_{i=1}^n (T_i / \hat{p}_{ti})} + \hat{\mu}_{t,S}.$$

REMARK 7. The estimator (4.7) resembles the familiar GREG estimator in survey sampling [Särndal, Swensson and Wretman (1992)], and it looks similar to the “doubly-robust” estimator (3.4). Notice, however, that (4.7) accounts for an informative assignment process as reflected by the use of the probabilities $\hat{p}_{ti} = \hat{\Pr}(T_i = t | y_{ti}, \mathbf{v}_i, i \in S)$ instead of the propensity scores $\hat{e}_{ti} = \hat{\Pr}(T_i = t | \mathbf{z}_i, i \in S)$. On the other hand, the estimator (4.7) does not possess a “double robustness” property, since the unknown model parameters are estimated jointly from the likelihood (4.3).

The estimators (4.5) and (4.6) are functions of the *mle* $\hat{\theta}_t$, and, hence, their large sample variance can be estimated by use of the inverse of the estimated information matrix. Large sample properties of the combined estimator (4.7) and a consistent estimator of its variance can be derived by application of *M*-estimation theory, see Landsman (2008) for details. The estimated variances of all the three estimators account for the sample distribution of the responses given the observed covariates.

4.4. *Application of new method to PISA data for Ireland.* We assume a normal distribution for the potential population responses under strongly ignorable assignment and a logistic model for the assignment probabilities:

$$(4.8) \quad \begin{aligned} f_p(y_{ti} | \mathbf{x}_i) &= N(\beta_{0t} + \mathbf{x}_i' \beta_t, \sigma_t^2); \\ \Pr(T_i = t | y_{ti}, \mathbf{v}_i, i \in S) &= \frac{\exp(\gamma_{0t} + \delta_t y_{ti} + \mathbf{v}_i' \gamma_t)}{1 + \exp(\gamma_{0t} + \delta_t y_{ti} + \mathbf{v}_i' \gamma_t)}, \quad t = 0, 1. \end{aligned}$$

The goodness of fit of the model is tested in Section 4.6.

When fitting the models to the two types of schools we found that the *x*-variables contain all the variables listed in Section 3.4 except for school location (S.loc), which was found to be highly insignificant in both the public and private school models. The *v*-variables contain all the variables listed in Section 3.4 except for mother’s education (ME), which was found to be highly insignificant in both the public and private school models. As discussed in Section 5, the sample *pdf* (4.1) obtained in the normal/logistic case is identifiable and accommodates consistent and asymptotically normal (CAN) estimators for all the model parameters, if *x* has at least one covariate not included in *v*.

Results. We computed the *mle* of the unknown parameters by maximizing the likelihood (4.3) with respect to $\theta_t = (\beta_{0t}, \beta_t, \sigma_t^2)$ and $\alpha_t = (\gamma_{0t}, \delta_t, \gamma_t)$, $t = 0, 1$. See Supplement B in the supplementary material [Pfeffermann and

TABLE 3
*Estimates and SE (in parenthesis) of model
coefficients*

Variable	Private schools	Public schools
Assignment (logistic) model		
Const	−2.95 (1.30)	13.88 (2.90)
δ_t	0.49 (0.21)	−2.02 (0.39)
GEN	0.77 (0.13)	−0.76 (0.18)
SEI	−0.12 (0.07)	0.40 (0.12)
HER	3.16 (0.20)	−2.57 (0.30)
SES	0.09 (0.07)	0.27 (0.11)
S.loc	1.13 (0.13)	−1.63 (0.24)
Population (normal) model		
σ_t	0.83 (0.02)	1.10 (0.07)
Const	6.09 (0.07)	6.89 (0.14)
GEN	−0.20 (0.05)	0.17 (0.08)
ME	0.18 (0.05)	0.11 (0.07)
SEI	0.16 (0.03)	0.16 (0.04)
HER	0.39 (0.09)	1.35 (0.20)
SES	0.21 (0.02)	0.30 (0.04)

Landsman (2011)] for the maximization procedure. Table 3 shows the estimates and standard errors (SE) of the model coefficients.

As anticipated, $\hat{\delta}_1 > 0$ and $\hat{\delta}_0 < 0$, but $\hat{\delta}_1$ is close to zero, although significant at the 5% level using the conventional t -statistic. On the other hand, $\hat{\delta}_0$ is highly negative, indicating that for given values of the covariates, the probability to attend a public school decreases very rapidly as the proficiency score increases. This finding suggests that pupils attending public schools are generally less able, and their lower scores are not necessarily because of poor quality of the public schools. Another important result emerging from Table 3 is that the variable school location (S.loc) is highly significant in the assignment models even when including the response among the covariates. We return to this result in Section 4.7. The coefficient of S.loc is positive for private schools and negative for public schools, suggesting that pupils from urban areas tend to enroll in private schools.

Table 4 shows the estimates of the population means by type of school and the corresponding estimates of the ATE. We present the two estimates defined in Section 4.3: the estimate $\hat{\mu}_{t,S}$ [equation (4.6)] and the combined estimate $\hat{\mu}_{t,C}$ [equation (4.7)].

The two methods of estimating the population means yield similar estimates and the ATE estimates are therefore likewise similar, negative and highly significant, indicating the very interesting and important result that after accounting for the school selection by pupils, the mean score in the pub-

TABLE 4
Estimation of population means and ATE

	Private school		Public school		ATE	
	$\hat{\mu}_{1,S}$	$\hat{\mu}_{1,C}$	$\hat{\mu}_{0,S}$	$\hat{\mu}_{0,C}$	$\hat{\tau}_S = \hat{\mu}_{1,S} - \hat{\mu}_{0,S}$	$\hat{\tau}_C = \hat{\mu}_{1,C} - \hat{\mu}_{0,C}$
Estimate	6.10	6.09	7.05	6.91	-0.95	-0.82
Std. error	0.05	0.06	0.15	0.12	0.16	0.13

lic schools is actually higher than in the private schools. A similar conclusion was reached in Section 3.4 by use of the LV and IV methods, although with smaller ATE estimates.

4.5. *Probability weighted estimation under the proposed approach.* We computed the PWE of the population means and the ATE under the proposed approach, accounting for possible sampling effects, similarly to the estimators computed under the previous methods shown in Section 3.5. See Supplement A of the supplementary material [Pfeffermann and Landsman (2011)] for the corresponding expressions. Table 5 shows the unweighted (UNWEI) and probability weighted estimators (WEI) of the school mean scores and the ATE. The results in this table reaffirm the results in Section 3.5 that the PISA sampling design is not informative for the models used for estimation of the ATE. We ignore the sample selection process in subsequent analysis.

4.6. *Goodness of fit of the model fitted to the PISA data for Ireland.* Assessing the goodness of fit of the model (4.8), or, more generally, any other model of the form (4.1) seems formidable on first sight since no observations are available from the population distribution under strong ignorability, and the assignment probabilities are generally unknown. Note, however, that once the identifiability of the sample *pdf* has been established (see Section 5), one faces the classical problem of having a random sample from a hypothesized *pdf* that needs to be tested. Below we consider three tests, which compare the theoretical and empirical sample distributions of the responses and apply them to the PISA data. The power of the tests is studied further in Appendix B by using simulated data sets.

TABLE 5
Unweighted and weighted estimators of schools mean score and ATE

Method	Private schools		Public schools		ATE	
	UNWEI	WEI	UNWEI	WEI	UNWEI	WEI
Est. pop. regression	6.10	6.09	7.05	6.92	-0.95	-0.83
Combined estimator	6.09	6.08	6.91	6.80	-0.82	-0.72

Goodness-of-fit tests. Suppose first that the true model parameters $\{\alpha_t, \theta_t\}$ are known. Denote by $U_{ti}(y) = F_{ti}(y|z_i) = \int_{-\infty}^y f_{S_t}(y_{ti}|z_i; \alpha_t, \theta_t) dy_{ti}$ the hypothesized *cdf* of $Y_{ti}|z_i, i = 1, 2, \dots, n_t$. For continuous F_{ti} , the random variables $U_{ti}(y)$ evaluated at the outcomes y_{ti} are independent Uniform[0, 1] variables since the responses Y_{ti} are independent given the covariates z_i . Let u_{t1}, \dots, u_{tn_t} denote the values of U_{t1}, \dots, U_{tn_t} computed at the sample values y_{t1}, \dots, y_{tn_t} and let $F_{t,\text{EMP}}$ define the empirical distribution of u_{t1}, \dots, u_{tn_t} . The following goodness-of-fit tests are in common use, where $u_{t(1)}, \dots, u_{t(n_t)}$ are the ordered values of u_{t1}, \dots, u_{tn_t} [Stephens (1986)]:

Kolmogorov–Smirnov:

$$(4.9) \quad KS_t = \max_i |F_{t,\text{EMP}}(u_{t(i)}) - u_{t(i)}|,$$

Cramer–von Misses:

$$(4.10) \quad CM_t = \frac{1}{12n_t} + \sum_{i=1}^{n_t} \left[u_{t(i)} - \frac{2i-1}{2n_t} \right]^2,$$

Anderson–Darling:

$$(4.11) \quad AD_t = -n_t - \frac{1}{n_t} \sum_{i=1}^{n_t} [(2i-1) \ln(u_{t(i)}) + (2n_t+1-2i) \ln(1-u_{t(i)})].$$

It is known [e.g., Babu and Feigelson (2006)] that KS is sensitive to large-scale differences in location and shape between the model and the empirical distribution, CM is sensitive to small-scale differences in the shape and AD is sensitive to differences near the tails of the distribution. All the three test statistics are *distribution-free* as long as the hypothetical *cdf* is fully specified (known parameters).

When computed with estimated parameters, the asymptotic distribution of the three statistics depends in a complex way on the true underlying *cdf*, and possibly also on the method of estimation. Correct critical values can be obtained in this case by use of parametric bootstrap. The procedure consists of generating a large number of samples from the estimated hypothesized model, re-estimating the unknown model parameters from each bootstrap sample and then computing the corresponding test statistics. The bootstrap distribution of the test statistics approximates the true distributions under the hypothesized model with correct order of error [Babu and Rao (2004)].

Validating the model fitted to the PISA data for Ireland. As explained above, the critical values for the distribution of the test statistics can be computed from the bootstrap distribution. To this end, we generated 250 bootstrap samples for each type of school ($t = 0, 1$) by first generating new outcomes Y_{ti} from the estimated normal population model $f_p(y_{ti}|x_i; \hat{\beta}_t, \hat{\sigma}_t^2)$,

TABLE 6
Goodness-of-fit test statistics and p-values (in parenthesis)

Private schools			Public schools		
KS	CM	AD	KS	CM	AD
0.023 (0.12)	0.089 (0.18)	0.62 (0.11)	0.027 (0.17)	0.062 (0.32)	0.45 (0.15)

using the same covariates as for the actual sample, and then selecting the units to the sample S_t using the estimated logistic probabilities $\Pr(T_i = t | v_i, y_{ti}; \hat{\delta}_t, \hat{\gamma}_t)$. Notice that this way the sample sizes in the two groups are no longer constant. We found that for the treatment group the mean sample size is 1,245 with standard deviation 17, and for the control group the mean sample size is 700 with standard deviation 18. (The sample sizes in the true data set are 1,244 and 694, resp.) Next we computed the *mle* of the parameters $\beta_t, \sigma_t^2, \gamma_t, \delta_t$ for each bootstrap sample and the test statistics (4.9)–(4.11).

Table 6 shows the values of the three test statistics for the PISA samples and their *p*-values, as computed from the corresponding bootstrap distributions. As can be seen, all the three statistics are nonsignificant with *p*-values higher than 10%, thus supporting the use of the selected models.

REMARK 8. We also computed the empirical means and standard deviations over the 250 bootstrap samples of the estimates of the model coefficients and all the ATE estimates considered in Sections 3 and 4. To save in space, we do not show the detailed results, but the means are generally close (and in most cases very close) to the values computed for the PISA data, and likewise for the standard deviations. These results indicate that the model coefficients can be estimated almost unbiasedly and with acceptable standard error estimates, despite the complicated structure of the sample model. Obtaining similar ATE estimates under the different methods to the estimates computed from the PISA data is another indication of the goodness of fit of the models fitted to this data set.

4.7. *Testing the assumptions underlying existing methods.* As mentioned earlier, the use of the proposed approach enables testing some of the assumptions underlying existing methods under the sample model fitted to the observed data. Consider first the logistic models for the assignment probabilities. The coefficients $\hat{\delta}_t$ of the response are significant in both models, with $\hat{\delta}_0$, in particular, being highly negative. This result suggests that the covariates available for the present study are not sufficient to explain the choice of school, and, hence, that methods that assume strong ignorability [assumptions SI(*a*)–SI(*b*)] and, in particular, methods that employ propensity scores computed with these covariates are not adequate. Notice in Ta-

ble 1 that the use of these methods yields positive ATE estimates, although of lower magnitude than the crude difference, $\bar{y}_1 - \bar{y}_0$.

Next consider the IV method. In its simple form it assumes the model $Y_t = \mathbf{x}'\beta_t + u_t$; $E_p(u_t|\mathbf{x}) = 0$, with three added conditions on the instrument h : IV(a)— $E_p(Y_t|\mathbf{x}, h) = E_p(Y_t|\mathbf{x})$, IV(b)— $E_p[T(u_1 - u_0)|\mathbf{x}, h] = 0$, and IV(c)— $\Pr(T = 1|\mathbf{x}, h) \neq \Pr(T = 1|\mathbf{x})$. The population model with the covariates \mathbf{x} listed in Section 3.4 is validated in Section 4.6. Furthermore, our analysis shows that the instrument S.loc is highly insignificant in the two population models, thus supporting the condition IV(a). Condition IV(b) cannot be verified empirically, but this condition is generally considered as being mild and it can be relaxed further [Wooldridge (2002)]. Finally, the condition IV(c) is verified as well since the coefficients of the instrument in the models fitted for the assignment probabilities are highly significant (Table 3), even when including the response y as an additional explanatory variable. Indeed, the use of the IV method yields an ATE estimate of -0.61 (Table 1), which is the closest to the estimate obtained under our approach among the other methods considered.

REMARK 9. We emphasize again that all the above conclusions are under the model that we have fitted (and validated) to the data.

5. Foundation of proposed approach.

5.1. *Identification of the sample distribution.* An important question underlying the use of the sample *pdf* (4.1) is model identification. In order to get some insight into this issue, we restrict to a given treatment t and hence drop for convenience the subscript t everywhere, denoting by $p(y, \mathbf{v}; \alpha) = \Pr(i \in S_t | y_t, \mathbf{v}; \alpha)$ the probability assignment rule (PAR) to the sample S_t , and by $f_p(y|\mathbf{x}; \theta)$ the corresponding population *pdf* of $Y_t|\mathbf{x}$ under strong ignorability. We assume that the response is continuous. Let $J \subseteq R$ define the common domain of the y -values for these functions. The sample *pdf* for units in S_t is therefore $f_{S_t}(y|\mathbf{x}, \mathbf{v}; \theta, \alpha) = \frac{f_p(y|\mathbf{x}; \theta)p(y, \mathbf{v}; \alpha)}{\int f_p(y|\mathbf{x}; \theta)p(y, \mathbf{v}; \alpha) dy}$ and the identifiability of the sample *pdf* is defined as follows:

DEFINITION 1. The sample *pdf* $f_{S_t}(y|\mathbf{x}, \mathbf{v}; \theta, \alpha)$ is identifiable if no different (proper) densities $f_p^{(1)}(y|\mathbf{x}; \theta^{(1)})$, $f_p^{(2)}(y|\mathbf{x}; \theta^{(2)})$ and different PARs $p^{(1)}(y, \mathbf{v}; \alpha^{(1)})$, $p^{(2)}(y, \mathbf{v}; \alpha^{(2)})$ exist such that the pairs $[f_p^{(1)}(y|\mathbf{x}; \theta^{(1)}), p^{(1)}(y, \mathbf{v}; \alpha^{(1)})]$ and $[f_p^{(2)}(y|\mathbf{x}; \theta^{(2)}), p^{(2)}(y, \mathbf{v}; \alpha^{(2)})]$ induce the same sample *pdf* for every $y \in J$ and every set of covariates (\mathbf{x}, \mathbf{v}) .

Clearly, if different pairs $[f_p^{(1)}, p^{(1)}], [f_p^{(2)}, p^{(2)}]$ yield the same sample *pdf*, the model is not identifiable. At first thought, this would seem to be always the case since (4.1) is the same *pdf* for the pair $[f_p^{(1)}(y|\mathbf{x}; \theta^{(1)}), p^{(1)}(y, \mathbf{v}; \alpha^{(1)})]$, and when the population *pdf* is $f_p^{(2)}(y|\mathbf{x}, \mathbf{v}; \theta^{(1)}, \alpha^{(1)}) =$

$\frac{f_p^{(1)}(y|x;\theta^{(1)})p^{(1)}(y,v;\alpha^{(1)})}{\int f_p^{(1)}(y|x;\theta^{(1)})p^{(1)}(y,v;\alpha^{(1)})dy}$ and the units are assigned with probabilities that do not depend on y given v (ignorable assignment). The *pdf* $f_p^{(2)}(y|x, v; \theta^{(1)}, \alpha^{(1)})$, however, does not generally belong to a conventional parametric family and is very different from $f_p^{(1)}(y|x; \theta^{(1)})$, especially when the assignment mechanism is strongly informative. Hence, field experts should be able to decide which of the two *pdfs*, $f_p^{(1)}(y|x; \theta^{(1)})$ or $f_p^{(2)}(y|x, v; \theta^{(1)}, \alpha^{(1)})$, is a more sensible population *pdf* for the potential responses in a given problem.

Conditions for model identification. Lemma 1 defines different conditions under which the sample *pdf* is identifiable. We assume for convenience that there are no covariates, but all the results generally hold when covariates x, v exist. Define $R_p(y; \theta^{(1)}, \theta^{(2)}) = \frac{f_p^{(2)}(y; \theta^{(2)})}{f_p^{(1)}(y; \theta^{(1)})}$. We assume throughout this section that the functions $f_p^{(j)}(y; \theta^{(j)})$ and $p^{(j)}(y; \alpha^{(j)})$, $j = 1, 2$, are strictly positive on $J^* \subseteq J$.

LEMMA 1. (a) Suppose that $J^* = [c, \infty)$ for some constant c . If $f_p^{(1)}(y; \theta^{(1)})$ and $f_p^{(2)}(y; \theta^{(2)})$ are two different *pdfs* satisfying for any given $\theta^{(1)}, \theta^{(2)}$, $\lim_{y \rightarrow \infty} R_p(y; \theta^{(1)}, \theta^{(2)}) = 0, \infty$ or does not exist, then there are no different PARs $p^{(1)}(y; \alpha^{(1)}), p^{(2)}(y; \alpha^{(2)})$ with finite positive limits as $y \rightarrow \infty$ yielding the same sample *pdf* for all $y \in J^*$.

(b) Suppose that $J^* = (-\infty, c]$ for some constant c . If $f_p^{(1)}(y; \theta^{(1)})$ and $f_p^{(2)}(y; \theta^{(2)})$ are two different *pdfs* satisfying for any given $\theta^{(1)}, \theta^{(2)}$, $\lim_{y \rightarrow -\infty} R_p(y; \theta^{(1)}, \theta^{(2)}) = 0, \infty$ or does not exist, then there are no different PARs $p^{(1)}(y; \alpha^{(1)}), p^{(2)}(y; \alpha^{(2)})$ with finite positive limits as $y \rightarrow -\infty$ yielding the same sample *pdf* for all $y \in J^*$.

(c) Let y_0 be a limit point of J^* . If $f_p^{(1)}(y; \theta^{(1)})$ and $f_p^{(2)}(y; \theta^{(2)})$ are two different *pdfs* satisfying for any given $\theta^{(1)}, \theta^{(2)}$, $\lim_{y \rightarrow y_0^+ (y \rightarrow y_0^-)} R_p(y; \theta^{(1)}, \theta^{(2)}) = 0, \infty$ or does not exist, then there are no different PARs $p^{(1)}(y; \alpha^{(1)}), p^{(2)}(y; \alpha^{(2)})$ with finite positive limits at $y = y_0$ yielding the same sample *pdf* for all $y \in J^*$.

PROOF. Part (a) is similar to Lee and Berger (2001) and is proved in Supplement C of the supplementary material [Pfeffermann and Landsman (2011)]. The proofs of the other two parts are similar. \square

Lemma 1 enables verifying the identifiability of the sample *pdf* for many combinations of population *pdfs* and PARs, but there are cases that need to be studied separately. For example, the lemma is not applicable to the case of normal population *pdfs* and logistic PARS if the coefficients of y in the two logistic distributions are allowed to have opposite signs. This is because

in this case one of the PARS will have a limit of zero as $y \rightarrow \infty$ or $y \rightarrow -\infty$, and the other PAR will have a limit of 1 (but see Result 1 below).

Further model identification results. Result 1 states the identifiability of the sample *pdf* resulting from the combination of a normal population *pdf* and a logistic *PAR*. The proof is given in Supplement D of the supplementary material [Pfeffermann and Landsman (2011)]. Landsman (2008) considers other combinations of population *pdfs* and *PARs*.

RESULT 1. No different pairs $[f_p(y|x;\theta), p(y, v; \alpha)]$ of a normal *pdf* and logistic *PAR* yield the same sample *pdf*, if the vectors x and v differ in at least one covariate.

REMARK 10. The condition on the covariates seems to impose a limitation on the model, but, in practice, there is no reason why the covariates used to model the response under strong ignorability should be the same as the covariates used to model the treatment assignment probabilities. See the models in Section 4.4.

5.2. *Practical identifiability.* Section 5.1 and the additional results in Landsman (2008) establish the “theoretical identifiability” of the sample model under a large number of plausible combinations of population *pdfs* and *PARs*. It is important to mention, however, that identifiability problems may arise in practice, depending on the forms of $f_p(y|x;\theta)$ and $p(y, v; \alpha)$. For example, Lee and Berger (2001) consider the case $f_p^{(1)}(y) = N(0, 1)$, $p^{(1)}(y) = \Phi(y - 1)$. The authors show graphically that in this case the sample density $f_S(y) = f_p^{(1)}(y)p^{(1)}(y) / \int f_p^{(1)}(y)p^{(1)}(y) dy$ can be closely approximated by the normal density $N(0.92, 0.79^2)$. This means that even though the sample *pdf* $f_S(y)$ is theoretically identifiable [Landsman (2008)], a problem may arise in practice when fitting the model, distinguishing between this density and the sample density obtained from $f_p^{(2)}(y) = N(0.92, 0.79^2)$, $p^{(2)}(y) = \text{const}$. Lee and Berger (2001) refer to this phenomenon as “practical non-identifiability.” Another example is where the *PAR* is logistic. Suppose that $p^{(1)}(y) = \{1 + [\exp(1 + y)]^{-1}\}^{-1}$, $p^{(2)}(y) = \{1 + [\exp(-1 - y)]^{-1}\}^{-1}$. Then, the pairs $[f_p^{(1)}(y) = N(0, 1), p^{(1)}(y)]$ and $[f_p^{(2)}(y) = N(1, 1), p^{(2)}(y)]$ induce the same sample *pdf* (“theoretical nonidentifiability”), and this *pdf* is closely approximated by $N(0.28, 0.93^2)$ (“practical nonidentifiability”).

We emphasize that in the presence of covariates, if the vectors x and v differ in at least one covariate, the problem of practical identifiability will generally not exist. See Landsman (2008) for a detailed analysis.

5.3. *Asymptotic properties of the maximum likelihood estimators.* The *mle* $\{\hat{\alpha}_t, \hat{\theta}_t, t = 0, 1\}$ obtained by maximization of (4.3) are the solutions of

the estimating equations $\sum_{k=1}^n u_{t,k} = 0$, where $u_{t,k}$ is the vector of the first derivatives of the k th component of the log-likelihood with respect to (α_t, θ_t) . Rotnitzky and Robins (1997) show that no \sqrt{n} -consistent and asymptotically normal (CAN) estimator of θ_t exists if (and only if) the derivatives of the log-likelihood with respect to θ_t are collinear with the derivatives of the log-likelihood with respect to α_t with probability 1. The derivatives are evaluated at the true parameter values. The authors illustrate that if the population model is $Y_{ti} \sim N(\beta_t, 1)$ and $\Pr(T_i = t|y_{ti}) = \frac{\exp(\gamma_{0t} + \delta_t y_{ti})}{1 + \exp(\gamma_{0t} + \delta_t y_{ti})}$, then no CAN estimator for β_t exists if in truth $\delta_t = 0$ (ignorable assignment).

Landsman (2008) shows that if $f_p(y_{ti}|x_i) = N(\beta_{0t} + x_i' \beta_t, \sigma_t^2)$, $\text{logit}[\Pr(T_i = t|y_{ti}, v_i)] = \gamma_{0t} + \delta_t y_{ti} + v_i' \gamma_t$ and x has at least one covariate not included in v , the derivatives of the log-likelihood (4.3) with respect to $\theta_t = (\beta_{0t}, \beta_t, \sigma_t)$ are not collinear with the derivatives with respect to $\alpha_t = (\gamma_{0t}, \gamma_t, \delta_t)$ with probability 1, and, hence, CAN estimators for the parameters exist even when the true assignment process is ignorable. This property enables testing the ignorability of the assignment process, as illustrated in Sections 4.4 and 4.7.

6. Summary remarks. In this article we study the use of an alternative approach for observational studies that recovers the treatment assignment model and the population model under strong ignorability from the sample data. It is shown that the sample model holding for the observed data, which incorporates the two models, is identifiable under some general conditions. Furthermore, the goodness of fit of the sample model can be tested successfully by standard test statistics because the sample model refers to the observed data. As illustrated in Section 4.7, the sample model enables also testing the appropriateness of the use of some of the existing methods for a particular data set.

We applied the new approach for comparing the proficiency in mathematics of children aged 15 between public and private schools in Ireland. Our analysis shows that although the average score of pupils in the sample from private schools is significantly higher than the average score of pupils from public schools, the picture is reversed once the effect of the school selection is accounted for properly. A similar conclusion is reached by the methods of latent variables and instrumental variables.

The approach proposed in this article is fully parametric, which raises questions of its robustness to departures from the models fitted to the data. We emphasize again that the models can be tested and, as the empirical illustrations show, the test statistics that we have applied have good power properties. Nonetheless, it is certainly worth considering the use of semi-parametric or nonparametric models either for the population models under strong ignorability and/or for the assignment probabilities, thus further robustifying the inference.

APPENDIX A: IMPUTATION OF PROFICIENCIES IN THE PISA STUDY

The multiple imputation approach draws at random multiple values from the conditional distribution of the unobserved proficiency, ψ_j of pupil j , given the m observed responses $d_j = (d_{1j}, \dots, d_{mj})$ and covariates \mathbf{x}_j representing individual background characteristics. The conditional *pdf* of ψ_j is expressed as

$$(A.1) \quad f(\psi_j | d_j, \mathbf{x}_j) \propto \prod_{i=1}^m [\Pr(D_{ij} = 1 | a_i, b_i, \psi_j)]^{d_{ij}} \times [\Pr(D_{ij} = 0 | a_i, b_i, \psi_j)]^{(1-d_{ij})} f(\psi_j | \mathbf{x}_j, \lambda, \sigma),$$

where $f(\psi_j | \mathbf{x}_j, \lambda, \sigma^2)$ is the normal distribution with mean $\mathbf{x}_j' \lambda$ and variance σ^2 , and $\Pr(D_{ij} = 1 | a_i, b_i, \psi_j) = [1 + \exp(-a_i(\psi_j - b_i))]^{-1}$. The parameter a_i measures how question i distinguishes between pupils and the parameter b_i represents the “difficulty” of question i . The responses to the various questions are assumed to be independent given (a_i, b_i, ψ_j) . Five imputed values of ψ_j are drawn for every student j in the sample and stored in the PISA database.

APPENDIX B: POWERS OF GOODNESS-OF-FIT TEST STATISTICS

In Section 4.6 we considered three goodness-of-fit test statistics and applied them for testing the model (4.8) fitted to the PISA data. In order to study the powers of the three test statistics in the case of misspecified population models, we simulated new data sets for each of the two groups (public and private schools) from the same models as in Section 4.6, except that the residual terms in the two population models were generated from the skew t -distribution defined by Azzalini and Capitanio (2003). The true population means and hence the ATE remain unchanged. The skew t -distribution depends on four parameters: location (ξ), scale (w), shape (α) and degrees of freedom (v). The normal and t -distributions are members of this family of distributions. For example, the case $\xi = 0$, $w = 1$, $\alpha = 0$, $v = \infty$ defines the standard normal distribution.

We generated 100 sets of residuals for each group from the following 3 distributions: (A) $\xi = 0$, $w = 1$, $\alpha = 0$, $v = 6$, (B) $\xi = -1.16$, $w = 1.55$, $\alpha = 2.5$, $v = \infty$, (C) $\xi = -1.24$, $w = 1.45$, $\alpha = 2.5$, $v = 6$. Distribution A defines the standard t -distribution with 6 degrees of freedom. Distribution B defines a skewed distribution with relatively short tails, while Distribution C defines a skewed distribution with a heavy right tail. The three densities are plotted in Figure 1. The location parameters were chosen in such a way that all the three distributions have mean zero, implying that the true ATEs are the same. The standard deviations equal 1.22, 1.04 and 1.27, respectively. Next we fitted the models that assume that the residuals are normal [equation (4.8)], such that the true models are misspecified.

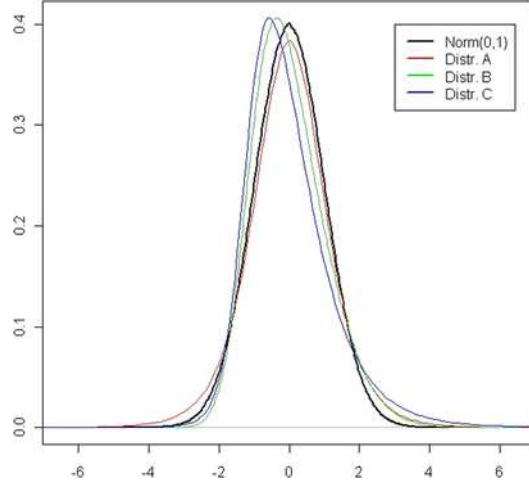


FIG. 1. *Densities of residuals under the four alternative distributions.*

Table B1 shows the percentage of samples that the misspecified models have been rejected by the three test statistics when using the conventional 10%, 5%, 2.5% and 1% significance levels. The percentages estimate the power of the tests. As can be seen, all the tests reject the three misspecified models for the private schools in literally all the samples, and the powers of the Cramer–von Mises test, and, in particular, the powers of the Anderson–Darling test are acceptable for the three misspecified models for the public schools as well. We mention in this regard that the three misspecified distributions were chosen to be sufficiently close to the normal distribution (see Figure 1). Further empirical results not shown indicate that for only mild larger distortions from the normal distribution, the powers of the three tests for the public schools are likewise close to 1.

Table B1 also shows the average of the estimates of the mean score in each group and the ATE, and the corresponding standard errors of the estimates (in parenthesis), over the 100 data sets. The empirical averages deviate now significantly from the corresponding true values ($\mu^1 = 6.10$, $\mu^0 = 7.05$; $\Delta = -0.95$) under all the three misspecified models and with larger standard errors than under the correct model. However, the estimates of the ATE are still highly negative, as under the correct model.

The conclusion from this simulation study is that even mild distortions from normality can affect the magnitude of the estimates of the ATE (but not their sign), but these mild distortions can be detected by the goodness-of-fit test statistics.

Acknowledgments. This paper is based on the Ph.D. thesis of the second author written at the Hebrew University of Jerusalem, Israel, under the

TABLE B1
Average estimates (SE) and percent of samples with rejected model

Private schools					Public schools			
$\hat{\mu}_{1,S} = 6.01 (0.14), \hat{\mu}_{1,C} = 5.99 (0.18)$					$\hat{\mu}_{0,S} = 7.34 (0.13), \hat{\mu}_{0,C} = 7.09 (0.39)$			
ATE $\hat{\tau}_S = -1.32 (0.20); \hat{\tau}_C = -1.10 (0.43)$								
Sig. level	0.10	0.05	0.025	0.01	0.10	0.05	0.025	0.01
Distribution A								
K-S	100	100	99	96	75	64	55	52
C-M	100	100	100	100	91	87	79	72
A-D	100	100	100	100	94	90	88	83
$\hat{\mu}_{1,S} = 5.96 (0.19), \hat{\mu}_{1,C} = 6.01 (0.15)$					$\hat{\mu}_{0,S} = 6.84 (0.10), \hat{\mu}_{0,C} = 6.89 (0.26)$			
ATE $\hat{\tau}_S = -0.88 (0.22); \hat{\tau}_C = -0.88 (0.31)$								
Sig. level	0.10	0.05	0.025	0.01	0.10	0.05	0.025	0.01
Distribution B								
K-S	100	100	100	99	72	64	51	47
C-M	100	100	100	100	83	77	66	63
A-D	100	100	100	100	92	79	76	70
$\hat{\mu}_{1,S} = 5.42 (0.20), \hat{\mu}_{1,C} = 5.49 (0.67)$					$\hat{\mu}_{0,S} = 6.76 (0.08), \hat{\mu}_{0,C} = 6.90 (0.46)$			
ATE $\hat{\tau}_S = -1.34 (0.22); \hat{\tau}_C = -1.41 (0.83)$								
Sig. level	0.10	0.05	0.025	0.01	0.10	0.05	0.025	0.01
Distribution C								
K-S	100	100	100	100	67	56	40	39
C-M	100	100	100	100	85	75	71	67
A-D	100	100	100	100	93	88	79	78

supervision of the first author. The authors are grateful to the Editor, Associate Editor and two referees for very insightful and constructive comments that improved the quality of the paper.

SUPPLEMENTARY MATERIAL

Supplement to: “Are private schools better than public schools? Appraisal for Ireland by methods for observational studies”
(DOI: [10.1214/11-AOAS456SUPP](https://doi.org/10.1214/11-AOAS456SUPP); .zip). This supplement contains a PDF which is divided into five sections:

- Supplement A develops the probability weighted estimators of the ATE.
- Supplement B describes the maximization of the likelihood (4.3).
- Supplement C contains the proof of Lemma 1.

Supplement D contains the proof of Result 1.

Supplement E describes the data file, which is provided.

The data file PISA_math2000.R contains the data.

REFERENCES

- ABADIE, A. and IMBENS, G. W. (2006). Large sample properties of matching estimators for average treatment effects. *Econometrica* **74** 235–267. [MR2194325](#)
- ADAMS, R. and WU, M., EDS. (2002). PISA 2000 Technical report, OECD, Paris.
- AZZALINI, A. and CAPITANIO, A. (2003). Distributions generated by perturbation of symmetry with emphasis on a multivariate skew t -distribution. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **65** 367–389. [MR1983753](#)
- BABU, G. J. and FEIGELSON, E. D. (2006). Astrostatistics: Goodness-of-fit and all that! In *Astronomical Data Analysis Software and Systems XV, ASP Conference Series* (C. GABRIEL, C. ARVISET, D. PONZ AND E. SOLANO, eds.) **351** 127–136. Astronomical Society of the Pacific, San Francisco.
- BABU, G. J. and RAO, C. R. (2004). Goodness-of-fit tests when parameters are estimated. *Sankhyā* **66** 63–74. [MR2082908](#)
- BREWER, K. R. W. (1963). Ratio estimation and finite populations: Some results deducible from the assumption of an underlying stochastic process. *Austral. J. Statist.* **5** 93–105. [MR0182078](#)
- DRONKERS, J. and AVRAM, S. (2010). A cross-national analysis of the relations of school choice and effectiveness differences between private-dependent and public schools. *Educational Research and Evaluation* **16** 151–175.
- GELMAN, A., CARLIN, J. B., STERN, H. S. and RUBIN, D. B. (2004). *Bayesian Data Analysis*, 2nd ed. Chapman & Hall/CRC, Boca Raton, FL. [MR2027492](#)
- GREENLEES, J. S., REECE, W. S. and ZIESCHANG, K. D. (1982). Imputation of missing values when the probability of response depends on the variable being imputed. *J. Amer. Statist. Assoc.* **77** 251–261.
- HAJEK, J. (1971). Comment on “An essay on the logical foundations of survey sampling, part one”. In *The Foundations of Survey Sampling* (V. P. Godambe and D. A. Sprott, eds.) 236. Holt, Rinehart and Winston, Toronto, ON.
- HANUSHEK, E. (2002). Publicly provided education. In *Handbook of Public Economics* (A. J. AUERBACH AND M. FELDSTEIN, eds.) 2045–2141. North-Holland, Amsterdam.
- HECKMAN, J. and VYTLACIL, E. (2006). Econometric evaluation of social programs. In *Handbook of Econometrics* **6B** (J. HECKMAN AND E. LEAMER, eds.) 4810–4861. North-Holland, Amsterdam.
- HOXBY, C. M. (2000). Does competition among public schools benefit students and taxpayers? *The American Economic Review* **90** 1209–1238.
- IMBENS, G. and ANGRIST, J. (1994). Identification and estimation of local average treatment effects. *Econometrica* **62** 467–475.
- LANDSMAN, V. (2008). Estimation of treatment effects in observational studies by recovering the assignment probabilities and the population model. Ph.D. dissertation, Hebrew Univ. Jerusalem, Israel.
- LEE, J. and BERGER, J. O. (2001). Semiparametric Bayesian analysis of selection models. *J. Amer. Statist. Assoc.* **96** 1397–1409. [MR1946585](#)
- LITTLE, R. J. (2004). To model or not to model? Competing modes of inference for finite population sampling. *J. Amer. Statist. Assoc.* **99** 546–556. [MR2109316](#)
- LUNCEFORD, J. K. and DAVIDIAN, M. (2004). Stratification and weighting via the propensity score in estimation of causal treatment effects: A comparative study. *Stat. Med.* **23** 2937–2960.

- MADDALA, G. S. (1983). *Limited-Dependent and Qualitative Variables in Econometrics. Econometric Society Monographs in Quantitative Economics* **3**. Cambridge Univ. Press, Cambridge. [MR0799154](#)
- MCCAFFREY, D., RIDGEWAY, G. and MORRAL, A. (2004). Propensity score estimation with boosted regression for evaluating causal effects in observational studies. *Psychological Methods* **9** 403–425.
- PFEFFERMANN, D. and LANDSMAN, V. (2011). Supplement to “Are private schools better than public schools? Appraisal for Ireland by methods for observational studies.” [DOI:10.1214/11-AOAS456SUPP](#).
- PFEFFERMANN, D. and SVERCHKOV, M. Y. (2003). Fitting generalized linear models under informative sampling. In *Analysis of Survey Data (Southampton, 1999)* 175–195. Wiley, Chichester. [MR1978851](#)
- PFEFFERMANN, D. and SVERCHKOV, M. (2009). Inference under informative sampling. In *Sample Surveys: Inference and Analysis. Handbook of Statistics* **29B** (D. PFEFFERMANN AND C. R. RAO, eds.) 455–487. North-Holland, Amsterdam.
- QIN, J. and ZHANG, B. (2007). Empirical-likelihood-based inference in missing response problems and its application in observational studies. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **69** 101–122. [MR2301502](#)
- R DEVELOPMENT CORE TEAM (2004). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- ROSENBAUM, P. R. (2002). *Observational Studies*, 2nd ed. Springer, New York. [MR1899138](#)
- ROSENBAUM, P. R. and RUBIN, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika* **70** 41–55. [MR0742974](#)
- ROTNITZKY, A. and ROBINS, J. (1997). Analysis of semi-parametric regression models with non-ignorable non-response. *Stat. Med.* **16** 81–102.
- RUBIN, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *J. Educational Psychology* **66** 688–701.
- SÄRNDAL, C.-E., SWENSSON, B. and WRETMAN, J. (1992). *Model Assisted Survey Sampling*. Springer, New York. [MR1140409](#)
- SMITH, T. M. F. and SUGDEN, R. A. (1988). Sampling and assignment mechanisms in experiments, surveys and observational studies. *Internat. Statist. Rev.* **56** 165–180.
- STATA CORP (2004). *Stata Statistical Software: Release 7*. StataCorp LP, College Station, TX.
- STEPHENS, M. A. (1986). Tests based on EDF statistics. In *Goodness-of-Fit Techniques* (R. B. D’AGOSTINO AND M. A. STEPHENS, eds.) 97–193. Dekker, New York.
- VANDENBERGHE, V. and ROBIN, S. (2004). Evaluating the effectiveness of private education across countries: A comparison of methods. *Labour Economics* **11** 487–506.
- WOOLDRIDGE, J. M. (2002). *Econometric Analysis of Cross Section and Panel Data*. MIT Press, Cambridge, MA.

HEBREW UNIVERSITY OF JERUSALEM
JERUSALEM, 91905
ISRAEL
AND
UNIVERSITY OF SOUTHAMPTON
SOUTHAMPTON, SO17 1BJ
UNITED KINGDOM
E-MAIL: msdanny@soton.ac.uk

NATIONAL CANCER INSTITUTE
NIH
ROCKVILLE, MARYLAND 20852
USA
E-MAIL: landsmnv@mail.nih.gov